

CSE- 1006 LAB Assignment 2.2

Academic year: 2021-2022 Semester: WIN

Faculty Name: Dr. Arun kumar Gopu Date: 17/3/2022

Student name: M.Taran Reg. no.: 19BCE7346

dFisher's Iris Dataset



IRIS FLOWERS

Description

This famous (Fisher's or Anderson's) iris data set gives the measurements in centimetres of the variables sepal length and width and petal length and width, respectively, for 50 flowers from each of 3 species of iris. The species are Iris setosa, versicolor, and virginica.

Format

iris is a data frame with 150 cases (rows) and 5 variables (columns) named Sepal.Length, Sepal.Width, Petal.Length, Petal.Width, and Species

EXERCISES

· Print the dataset iris



```
> print(iris)
   Sepal.Length Sepal.Width Petal.Length
          5.1 3.5
1
2
          4.9
                    3.0
                              1.4
3
          4.7
                    3.2
                              1.3
          4.6
4
                    3.1
                              1.5
5
          5.0
                    3.6
                              1.4
                                     > iris
                    3.9
6
          5.4
                              1.7
                                         Sepal.Length Sepal.Width Petal.Length
7
                    3.4
                                                 5.1
          4.6
                              1.4
                                                            3.5
8
          5.0
                    3.4
                              1.5
                                     2
                                                 4.9
                                                            3.0
                                                                        1.4
9
          4.4
                    2.9
                              1.4
                                     3
                                                 4.7
                                                            3.2
                                                                        1.3
10
          4.9
                    3.1
                              1.5
                                                 4.6
                                                                        1.5
                                    4
                                                            3.1
11
          5.4
                    3.7
                              1.5
                                  5
                                                 5.0
                                                            3.6
12
          4.8
                              1.6
                    3.4
                                                            3.9
                                                                        1.7
                                  6
                                                 5.4
13
          4.8
                   3.0
                              1.4
                                    7
                                                            3.4
                                                 4.6
                                                                        1.4
14
          4.3
                   3.0
                              1.1 8
                                                 5.0
                                                            3.4
                                                                        1.5
15
          5.8
                   4.0
                              1.2 9
                                                            2.9
                                                 4.4
                              1.5 10
1.3 11
                   4.4
16
         5.7
                                                 4.9
                                                            3.1
                                                                        1.5
                    3.9
17
         5.4
                                                 5.4
                                                            3.7
                                                                        1.5
18
         5.1
                    3.5
                              1.4
                                     12
                                                 4.8
                                                            3.4
                                                                        1.6
19
         5.7
                    3.8
                              1.7
                                     13
                                                 4.8
                                                            3.0
                                                                        1.4
20
         5.1
                    3.8
                              1.5
                                     14
                                                 4.3
                                                           3.0
                                                                        1.1
                             1.7
21
         5.4
                    3.4
                                     15
                                                 5.8
                                                            4.0
                                                                        1.2
22
         5.1
                    3.7
                             1.5
                                    16
                                                 5.7
                                                            4.4
                                                                        1.5
23
         4.6
                    3.6
                             1.0
                                    17
                                                 5.4
                                                            3.9
                                                                        1.3
         5.1
                    3.3
24
                              1.7
                                     18
                                                 5.1
                                                            3.5
                                                                       1.4
25
                    3.4
          4.8
                              1.9
                                     19
                                                 5.7
                                                            3.8
                                                                        1.7
```

· Print the structure of the dataset iris

```
> str(iris)
'data.frame': 150 obs. of 5 variables:
$ sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6
5 4.4 4.9 ...
$ sepal.width: num 3.5 3 3.2 3.1 3.6 3.9 3.4
3.4 2.9 3.1 ...
$ Petal.Length: num 1.4 1.4 1.3 1.5 1.4 1.7 1.
4 1.5 1.4 1.5 ...
$ Petal.width: num 0.2 0.2 0.2 0.2 0.2 0.4 0.
3 0.2 0.2 0.1 ...
$ species : Factor w/ 3 levels "setosa","ve rsicolor",..: 1 1 1 1 1 1 1 1 1 1 ...
>
```

· Print the summary of all the variables of the dataset iris (Hint: Use function summary())



```
> summary(iris)
  Sepal.Length
                 Sepal.Width
      :4.300
                 Min. :2.000
Min.
1st Qu.:5.100
                 1st Qu.: 2.800
Median :5.800
                 Median :3.000
        :5.843
                        :3.057
Mean
                 Mean
 3rd Qu.:6.400
                 3rd Qu.:3.300
мах.
       :7.900
                Max.
                        :4.400
 Petal.Length
                 Petal.Width
Min.
       :1.000
                 Min.
                        :0.100
1st Qu.:1.600
                 1st Qu.: 0.300
Median :4.350
                 Median :1.300
Mean
        :3.758
                 Mean
                        :1.199
 3rd Qu.:5.100
                 3rd Qu.:1.800
       :6.900
                 Max.
                        :2.500
       Species
setosa
          :50
 versicolor:50
virginica:50
```

· How many of the variables (columns) are in the dataset iris

```
> ncol(iris)
[1] 5
```

· How many observations (rows) are in the dataset iris

```
> nrow(iris)
[1] 150
```

> duplicated(iris)

· Use duplicated() function to print the logical vector indicating the duplicate values present in the dataset iris

```
[1] FALSE FALSE FALSE FALSE FALSE
[7] FALSE FALSE FALSE FALSE FALSE
[13] FALSE FALSE FALSE FALSE FALSE
[19] FALSE FALSE FALSE FALSE FALSE
[25] FALSE FALSE FALSE FALSE FALSE
[31] FALSE FALSE FALSE FALSE FALSE
[37] FALSE FALSE FALSE FALSE FALSE
[43] FALSE FALSE FALSE FALSE FALSE
[49] FALSE FALSE FALSE FALSE FALSE
[55] FALSE FALSE FALSE FALSE FALSE
[61] FALSE FALSE FALSE FALSE FALSE
[67] FALSE FALSE FALSE FALSE FALSE
[73] FALSE FALSE FALSE FALSE FALSE
[79] FALSE FALSE FALSE FALSE FALSE
[85] FALSE FALSE FALSE FALSE FALSE
[91] FALSE FALSE FALSE FALSE FALSE
```

[97] FALSE FALSE FALSE FALSE FALSE
[103] FALSE FALSE FALSE FALSE FALSE

· Extract duplicate elements from the dataset iris

```
> iris[duplicated(iris),]
    Sepal.Length Sepal.Width Petal.Length
143           5.8           2.7      5.1
    Petal.Width Species
143           1.9 virginica
>
```

· Extract unique elements from the dataset iris

```
> iris[!duplicated(iris),]
    Sepal.Length Sepal.Width Petal.Length
              5.1
                          3.5
2
             4.9
                          3.0
                                        1.4
3
             4.7
                          3.2
                                        1.3
                                                Petal.Width
                                                                Species
4
             4.6
                          3.1
                                        1.51
                                                        0.2
                                                                 setosa
5
             5.0
                          3.6
                                        1.4 2
                                                        0.2
                                                                 setosa
6
             5.4
                          3.9
                                        1.73
                                                        0.2
                                                                 setosa
7
             4.6
                          3.4
                                        1.44
                                                        0.2
                                                                 setosa
8
             5.0
                          3.4
                                        1.55
                                                        0.2
                                                                 setosa
9
             4.4
                          2.9
                                        1.46
                                                        0.4
                                                                 setosa
10
             4.9
                          3.1
                                        1.57
                                                        0.3
                                                                 setosa
11
             5.4
                          3.7
                                        1.58
                                                        0.2
                                                                 setosa
12
             4.8
                          3.4
                                        1.69
                                                        0.2
                                                                 setosa
                                        1.4.10
13
             4.8
                          3.0
                                                        0.1
                                                                 setosa
14
             4.3
                          3.0
                                        1.1 11
                                                        0.2
                                                                 setosa
15
             5.8
                          4.0
                                        1.2 12
                                                        0.2
                                                                 setosa
                                       1.5 13
16
             5.7
                          4.4
                                                        0.1
                                                                 setosa
17
             5.4
                          3.9
                                       1.3 14
                                                        0.1
                                                                 setosa
                                       1.4 15
18
             5.1
                          3.5
                                                        0.2
                                                                 setosa
19
             5.7
                          3.8
                                       1.7 16
                                                        0.4
                                                                 setosa
20
             5.1
                          3.8
                                        1.5 17
                                                        0.4
                                                                 setosa
21
             5.4
                          3.4
                                        1.7 18
                                                        0.3
                                                                setosa
22
             5.1
                          3.7
                                        1.5 19
                                                        0.3
                                                                setosa
23
             4.6
                          3.6
                                        1.0 20
                                                        0.3
                                                                 setosa
```

· Print the indices of duplicate elements in the dataset iris

```
> which(duplicated(iris))
[1] 143
```

· Print the indices of unique elements in the dataset iris



```
> which(!duplicated(iris))
  [1]
        1
             2
                  3
                      4
                           5
                                6
                                         8
 [10]
       10
            11
                 12
                     13
                          14
                               15
                                   16
                                        17
                                             18
 [19]
       19
            20
                 21
                     22
                          23
                               24
                                   25
                                        26
                                             27
 [28]
       28
            29
                 30
                     31
                          32
                               33
                                        35
                                             36
 [37]
        37
            38
                 39
                     40
                          41
                               42
                                   43
                                        44
                                             45
 [46]
       46
            47
                 48
                     49
                          50
                               51
                                   52
                                        53
 [55]
        55
            56
                 57
                      58
                          59
                               60
                                   61
                                        62
                                             63
 [64]
        64
            65
                 66
                     67
                          68
                               69
                                   70
                                        71
                                             72
 [73]
        73
            74
                 75
                     76
                               78
                                        80
                          77
                                             81
 [82]
        82
            83
                 84
                     85
                          86
                               87
 [91]
        91
            92
                 93
                     94
                          95
                               96
                                        98
                                   97
[100] 100 101 102 103 104 105 106 107 108
[109] 109 110 111 112 113 114 115 116 117
[118] 118 119 120 121 122 123 124 125 126
      127 128 129 130 131 132 133 134 135
[136] 136 137 138 139 140 141 142 144 145
[145] 146 147 148 149 150
```

· How many unique elements are in the dataset iris

```
> sum(!duplicated(iris))
[1] 149
```

· How many duplicate elements are in the dataset iris

```
> sum(duplicated(iris))
[1] 1
>
```

Missing Values:

- · A missing value is one whose value is unknown.
- · Missing values are represented in R by the NA symbol.
- · NA is a special value whose properties are different from other values.
- · NA is one of the very few reserved words in R: you cannot give anything this name.
- · Missing values are often legitimate: values really are missing in real life.
- · NAs can arise when you read in an Excel spreadsheet with empty cells, for example.
- · You will also see NA when you try certain operations that are illegal or don't make sense.

Here are some examples of operations that produce NA's.

EXERCISES

· Practice above examples that generate NA values

```
VIT-AP UNIVERSITY
```

```
> x <- c(12, 34, NA, 3, 4, NA, 56, NA,37,89,NA,43)
>
> is.na(x)
[1] FALSE FALSE TRUE FALSE FALSE TRUE
[7] FALSE TRUE FALSE FALSE TRUE FALSE
```

· Create NA values by some illegal operations

```
> as.numeric (c("2", "6", "three", "4"))
[1] 2 6 NA 4
Warning message:
NAs introduced by coercion
```

· Practice exercises in lecture slide

```
> X
[1] 12 34 NA 3 4 NA 56 NA 37 89 NA 43
> x + 1
[1] 13 35 NA 4 5 NA 57 NA 38 90 NA 44
>
> sum(x)
[1] NA
>
> length(x)
[1] 12
>
> is.na(x)
[1] FALSE FALSE TRUE FALSE FALSE TRUE FALSE
[8] TRUE FALSE FALSE TRUE FALSE
> which(is.na(x))
[1] 3 6 8 11
>
> x[! is.na(x)]
[1] 12 34 3 4 56 37 89 43
>
```

· What happens when we try to sort the data with NA values

Sorting data containing missing values in R is again different from other packages because NA cannot be compared to other values.

By default, sort removes any NA values and can therefore change the length of a vector.

```
> (temp <- sort(x))
[1] 3 4 12 34 37 43 56 89
```

· How to find the length of a vector with NA values

```
> length(temp)
[1] 8
```



The user can specify if NA should be last or first in a sorted order by indicating TRUE or FALSE for the na.last argument.

```
> sort(x, na.last = TRUE)
[1] 3 4 12 34 37 43 56 89 NA NA NA NA
```

DATASET - INTRODUCTION

In today's lab we are going to work on dataset "airquality"

Details of Dataset:

Daily readings of the following air quality values for May 1, 1973 (a Tuesday) to September 30, 1973.

Description of variables:

Ozone: Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island

Solar.R: Solar radiation in Langleys in the frequency band 4000--7700 Angstroms from 0800 to 1200 hours at Central Park

Wind: Average wind speed in miles per hour at 0700 and 1000 hours at LaGuardia Airport

Temp: Maximum daily temperature in degrees Fahrenheit at LaGuardia Airport.

EXERCISES

· Print the dataset airquality



```
> airquality
    Ozone Solar.R Wind Temp Month Day
        41
                190 7.4
                             67
2
        36
                118 8.0
                             72
                                     5
                                         2
3
        12
                149 12.6
                             74
                                     5
                                         3
4
        18
                313 11.5
                             62
                                     5
                                         4
5
        NA
                 NA 14.3
                             56
                                     5
                                         5
                                                             230 10.9
                                                                          75
                                                                                      9
                                             132
                                                     21
6
        28
                 NA 14.9
                             66
                                     5
                                         6
                                                             259 9.7
                                                     24
                                                                                  9
                                                                                     10
                                             133
                                                                          73
7
        23
                299
                     8.6
                             65
                                     5
                                         7
                                                     44
                                                             236 14.9
                                                                                  9
                                                                                     11
                                             134
                                                                          81
8
        19
                 99 13.8
                             59
                                     5
                                         8
                                             135
                                                     21
                                                             259 15.5
                                                                          76
                                                                                  9
                                                                                     12
9
         8
                 19 20.1
                             61
                                     5
                                         9
                                            136
                                                     28
                                                             238
                                                                  6.3
                                                                          77
                                                                                  Q.
                                                                                     13
10
        NA
                194
                     8.6
                             69
                                     5
                                        10
                                             137
                                                     9
                                                             24 10.9
                                                                          71
                                                                                     14
11
         7
                 NΑ
                     6.9
                             74
                                     5
                                        11
                                                             112 11.5
                                                                                  9
                                            138
                                                     13
                                                                          71
                                                                                     15
12
        16
                256
                     9.7
                             69
                                     5
                                        12
                                             139
                                                                  6.9
                                                                                  9
                                                     46
                                                             237
                                                                          78
                                                                                     16
13
        11
                290 9.2
                             66
                                     5
                                        13
                                             140
                                                     18
                                                             224 13.8
                                                                          67
                                                                                  9
                                                                                     17
                                                              27 10.3
14
        14
                274 10.9
                             68
                                     5
                                        14
                                             141
                                                     13
                                                                          76
                                                                                  9
                                                                                     18
                                                             238 10.3
                                            142
                                                     24
                                                                          68
                                                                                  9
                                                                                     19
15
        18
                 65 13.2
                             58
                                     5
                                        15
                                            143
                                                    16
                                                             201
                                                                 8.0
                                                                          82
                                                                                     20
16
        14
                334 11.5
                             64
                                     5
                                        16
                                             144
                                                     13
                                                             238 12.6
                                                                          64
                                                                                  9
                                                                                     21
17
        34
                307 12.0
                             66
                                     5
                                        17
                                             145
                                                                          71
                                                                                  9
                                                     23
                                                              14 9.2
                                                                                     22
18
         6
                 78 18.4
                             57
                                     5
                                        18
                                             146
                                                             139 10.3
                                                                                  9
                                                                                     23
                                                     36
                                                                          81
                322 11.5
19
        30
                             68
                                     5
                                        19
                                             147
                                                     7
                                                              49 10.3
                                                                          69
                                                                                  9
                                                                                     24
20
        11
                 44
                     9.7
                             62
                                     5
                                        20
                                            148
                                                     14
                                                              20 16.6
                                                                          63
                                                                                  9
                                                                                     25
21
         1
                  8
                     9.7
                             59
                                     5
                                        21
                                            149
                                                     30
                                                             193 6.9
                                                                          70
                                                                                     26
22
        11
                320 16.6
                            73
                                     5
                                        22
                                                                                  9
                                             150
                                                     NΑ
                                                             145 13.2
                                                                          77
                                                                                     27
23
         4
                 25 9.7
                             61
                                     5
                                        23
                                                     14
                                                                          75
                                                                                  9
                                             151
                                                             191 14.3
                                                                                     28
24
        32
                 92 12.0
                             61
                                     5
                                        24
                                             152
                                                     18
                                                             131
                                                                  8.0
                                                                          76
                                                                                  9
                                                                                     29
25
        NA
                 66 16.6
                             57
                                     5
                                        25
                                             153
                                                     20
                                                             223 11.5
                                                                          68
                                                                                  9
                                                                                     30
```

· Print the structure of the dataset airquality

```
> str(airquality)
               153 obs. of 6 variables:
'data.frame':
$ Ozone : int 41 36 12 18 NA 28 23 19 8 NA
$ Solar.R: int 190 118 149 313 NA NA 299 99
19 194 ...
        : num 7.4 8 12.6 11.5 14.3 14.9 8.6
$ Wind
13.8 20.1 8.6 ...
        : int 67 72 74 62 56 66 65 59 61 69
$ Temp
                5 5 5 5 5 5 5 5 5 5 ...
$ Month : int
$ Day
                1 2 3 4 5 6 7 8 9 10 ...
          : int
```

· Print the summary of all the variables of the dataset airquality (Hint: Use function summary())



```
> summary(airquality)
                     Solar.R
     Ozone
Min.
      : 1.00
                 Min. : 7.0
 1st Qu.: 18.00
                 1st Qu.:115.8
Median : 31.50
                 Median :205.0
      : 42.13
                       :185.9
Mean
                 Mean
 3rd Qu.: 63.25
                  3rd Qu.:258.8
       :168.00
                       :334.0
мах.
                 мах.
 NA's
       :37
                  NA's
                         :7
      Wind
                       Temp
        : 1.700
Min.
                 Min.
                         :56.00
 1st Qu.: 7.400
                 1st Qu.:72.00
Median : 9.700
                 Median :79.00
      : 9.958
                 Mean
                       :77.88
Mean
 3rd Qu.:11.500
                 3rd Qu.:85.00
       :20.700
                        :97.00
мах.
                 мах.
     Month
                      Day
                       : 1.0
Min.
       :5.000
                Min.
 1st Qu.:6.000
                 1st Qu.: 8.0
Median :7.000
                Median:16.0
      :6.993
Mean
                Mean
                      :15.8
 3rd Qu.:8.000
                3rd Qu.:23.0
```

· How many of the variables (columns) are in the dataset airquality

:31.0

мах.

```
> ncol(airquality)
[1] 6
```

:9.000

Max.

· How many observations (rows) are in the dataset airquality

```
> nrow(airquality)
[1] 153
```

· Use the function is.na() to find whether any missing values are in the dataset airquality



```
> is.na(airquality)
      Ozone Solar.R Wind Temp Month
  [1,] FALSE
            FALSE FALSE FALSE
  [2,] FALSE
             FALSE FALSE FALSE
  [3,] FALSE
             FALSE FALSE FALSE
  [4,] FALSE
            FALSE FALSE FALSE
 [5,] TRUE
[6,] FALSE
              TRUE FALSE FALSE FALSE
              TRUE FALSE FALSE FALSE
 [7,] FALSE
[8,] FALSE
             FALSE FALSE FALSE
            FALSE FALSE FALSE
  [9,] FALSE
             FALSE FALSE FALSE
 [10,] TRUE
            FALSE FALSE FALSE
 [11,] FALSE
[12,] FALSE
             TRUE FALSE FALSE FALSE
            FALSE FALSE FALSE
 [13,] FALSE
            FALSE FALSE FALSE
 [14,] FALSE
[15,] FALSE
[16,] FALSE
            FALSE FALSE FALSE
            FALSE FALSE FALSE
            FALSE FALSE FALSE
 [17,] FALSE
[18,] FALSE
            FALSE FALSE FALSE
            FALSE FALSE FALSE
 [19,] FALSE
             FALSE FALSE FALSE
 [20,] FALSE
             FALSE FALSE FALSE
 [21,] FALSE
            FALSE FALSE FALSE
 [22,] FALSE
             FALSE FALSE FALSE
 [23,] FALSE
             FALSE FALSE FALSE
```

· Print the indices of the missing values in the dataset airquality in column major representation

```
> colSums(is.na(airquality))
 Ozone Solar.R
                   Wind
                                  Month
                           Temp
    37
                              0
                                      0
                      0
   Day
> which(is.na(airquality))
 [1]
       5 10 25 26 27
                          32
                              33
                                  34
[11]
      37
         39 42
                 43
                     45
                         46
                                  53
                              52
                                           55
     56 57 58
                 59 60
                                  72
                                      75
[21]
                          61
                              65
[31] 84 102 103 107 115 119 150 158 159 164
[41] 180 249 250 251
```

· Print the indices of the missing values in the dataset airquality in row major representation



```
> which(is.na(airquality$0zone))
 [1]
       5 10
              25
                   26
                      27
                            32
                                33
                                    34
[11]
      37
          39
              42
                   43
                       45
                            46
                                52
                                     53
                                         54
                                             55
[21]
      56
          57
               58
                   59
                       60
                            61
                                65
                                         75
                                             83
      84 102 103 107 115 119 150
```

· Print indices of the missing values in row and column number wise (Hint: Use function which() and argument arr.ind = TRUE)

```
> which(is.na(airquality),arr.ind = TRUE)
       row col
 [1,]
         5
              1
 [2,]
              1
        10
 [3,]
        25
              1
 [4,]
        26
              1
 [5,]
        27
              1
 [6,]
        32
              1
        33
              1
 [8,]
        34
              1
 [9,]
        35
              1
[10,]
        36
              1
[11,]
        37
              1
[12,]
        39
              1
[13,]
        42
              1
[14,]
        43
              1
[15,]
        45
              1
[16,]
        46
              1
[17,]
        52
              1
[18,]
        53
              1
[19,]
        54
              1
[20,]
        55
              1
[21,]
        56
              1
[22,]
```

· How many missing values are in the dataset airquality?

```
> sum(is.na(airquality))
[1] 44
```

· Which variables are the missing values concentrated in?

```
> which(is.na(airquality))
[1]
      5 10 25
                26
                         32
                                     35
                                         36
                    27
                             33
                                 34
[11]
     37
         39 42
                43
                    45 46
                            52
                                 53
                                    54
                                         55
[21]
     56
        57
            58
                59 60 61
                             65
                                 72
                                     75
                                         83
     84 102 103 107 115 119 150 158 159 164
[41] 180 249 250 251
```

· How would you omit all rows containing missing values?



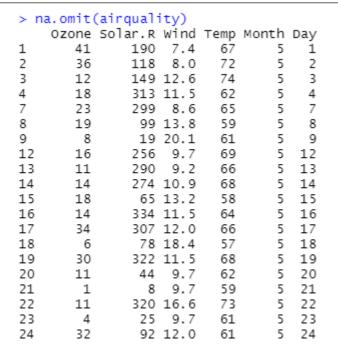
```
> na.omit(airquality)
    Ozone Solar.R Wind Temp Month Day
                190 7.4
                                     5
1
        41
                            67
2
        36
                118 8.0
                            72
                                     5
                                         2
3
        12
                149 12.6
                            74
                                     5
                                         3
4
        18
                313 11.5
                            62
                                     5
                                         4
                299 8.6
                                     5
                                         7
        23
                            65
8
        19
                 99 13.8
                                     5
                                         8
                             59
9
                 19 20.1
                                     5
                                         9
        8
                            61
                256 9.7
                                     5
                                        12
12
        16
                            69
                290 9.2
                                     5
13
        11
                                        13
                            66
                274 10.9
                                     5
14
        14
                            68
                                        14
                 65 13.2
                                     5
15
        18
                                        15
                             58
                334 11.5
                                     5
16
                                        16
        14
                             64
                307 12.0
                                     5
17
        34
                                        17
                            66
                                     5
18
                 78 18.4
                                        18
        6
                             57
                                     5
19
                322 11.5
                                        19
        30
                             68
                                     5
        11
                                        20
20
                 44
                     9.7
                            62
                                     5
                    9.7
21
        1
                  8
                             59
                                        21
                                     5
                320 16.6
22
        11
                            73
                                        22
                                     5
                 25 9.7
23
         4
                                        23
                             61
                                     5
                 92 12.0
24
        32
                             61
                                        24
                                     5
28
        23
                 13 12.0
                             67
                                        28
                                     5
        45
                252 14.9
                                        29
29
                             81
30
       115
                223 5.7
                             79
                                        30
```

· Print the records without missing values in the dataset airquality using the function complete.cases()

```
> airquality[complete.cases(airquality),]
    Ozone Solar.R Wind Temp Month Day
               190 7.4
1
        41
                            67
                                    5
                                         1
2
        36
                118 8.0
                            72
                                    5
                                         2
3
        12
                149 12.6
                            74
                                    5
                                         3
4
       18
                313 11.5
                            62
                                    5
                                        4
        23
                299 8.6
                            65
                                    5
                                         7
8
       19
                 99 13.8
                            59
                                    5
                                         8
9
        8
                19 20.1
                            61
                                    5
                                        9
12
       16
                256 9.7
                            69
                                    5
                                       12
13
       11
                290 9.2
                                    5
                                       13
                            66
14
       14
                274 10.9
                            68
                                    5
                                       14
15
       18
                65 13.2
                            58
                                    5
                                       15
16
       14
                334 11.5
                            64
                                    5
                                       16
17
        34
                307 12.0
                            66
                                    5
                                       17
                78 18.4
18
        6
                            57
                                    5
                                       18
                322 11.5
19
        30
                                    5
                                       19
                            68
20
        11
                 44
                    9.7
                                    5
                                       20
                            62
21
                  8
                     9.7
                            59
                                    5
                                       21
```

· Print the records without missing values in the dataset airquality using the function na.omit()





· Print the records without missing values in the dataset airquality using the function na.exclude()

> na.exclude(airquality)

	Ozone	Solar.R	Wind	Temp	Month	Day
1	41	190	7.4	67	5	1
2	36	118	8.0	72	5	2
3	12	149	12.6	74	5	3
4	18	313	11.5	62	5	4
7	23	299	8.6	65	5	7
8	19	99	13.8	59	5	8
9	8	19	20.1	61	5	9
12	16	256	9.7	69	5	12
13	11	290	9.2	66	5	13
14	14	274	10.9	68	5	14
15	18	65	13.2	58	5	15
16	14	334	11.5	64	5	16
17	34	307	12.0	66	5	17
18	6	78	18.4	57	5	18
19	30	322	11.5	68	5	19
20	11	44	9.7	62	5	20
21	1	8	9.7	59	5	21
22	11	320	16.6	73	5	22
23	4	25	9.7	61	5	23
24	32	92	12.0	61	5	24
28	23	13	12.0	67	5	28

· Print the records containing missing values in the dataset airquality using the function complete.cases()



> a	irqual	lity[!compl	ete.	cases((airqual	lity)	,]

_		, [- C		
	Ozone	Solar.R	Wind			Day	
5	NA	NA	14.3	56	5	5	
6	28		14.9	66	5	6	
10	NA	194		69	5	10	
11	7	NA	6.9	74	5	11	
25	NA	66	16.6	57	5	25	
26	NA	266	14.9	58	5	26	
27	NA	NA	8.0	57	5	27	
32	NA	286	8.6	78	6	1	
33	NA	287	9.7	74	6	2	
34	NA	242		67	6	3	
35	NA	186		84	6	4	
36	NA	220	8.6	85	6	5	
37	NA		14.3	79	6	6	
39	NA	273		87	6	8	
42	NA	259	10.9	93	6	11	
43	NA	250	9.2	92	6	12	
45	NA	332	13.8	80	6	14	
46	NA	322	11.5	79	6	15	
52	NA	150		77	6	21	
53	NA	59	1.7	76	6	22	
54	NA	91	4.6	76	6	23	
55	NA	250	6.3	76	6	24	
56	NA	135	8.0	75	6	25	
57	NA	127	8.0	78	6	26	
58	NA	47	10.3	73	6	27	
59	NA	98	11.5	80	6	28	
60	NA	31	14.9	77	6	29	
61	NA	138	8.0	83	6	30	
65	NA	101	10.9	84	7	4	
72	NA	139	8.6	82	7	11	
75	NA	291	14.9	91	7	14	
83	NA	258	9.7	81	7	22	
84	NA	295	11.5	82	7	23	
96	78	NA	6.9	86	8	4	
97	35	NA	7.4	85	8	5	
98	66	NA	4.6	87	8	6	
102	NA		8.6	92	8	10	
103	NA	137		86	8	11	
107	NA	64	11.5	79	8	15	
115	NA	255	12.6	75	8	23	
119	NA	153	5.7	88	8	27	
150	NA		13.2	77	9	27	